

What is claimed is:

- 1 1. A method of reading data comprising the steps of:
2 receiving a request for a data block stored within a stripe of erasure
3 coded data, the stripe of erasure coded data stored across a plurality of
4 storage devices which include a target storage device that holds the data
5 block;
6 sending read messages to the storage devices; and
7 receiving a reply message from each of at least a quorum of the storage
8 devices, the reply message from the target storage device including the
9 data block, the quorum meeting a quorum condition of a number such that
10 any two selections of the number of stripe blocks intersect in a minimum
11 number of the stripe blocks needed to decode the stripe.
- 1 2. The method of claim 1 wherein the stripe blocks comprise a first number of
2 data blocks and a second number of parity blocks.
- 1 3. The method of claim 1 wherein the read message that is sent to the target
2 storage device includes an indicator.
- 1 4. The method of claim 1 wherein each of the reply messages of the quorum
2 indicate that there is no pending write for the stripe block stored on the storage
3 device.
- 1 5. The method of claim 1 wherein the reply messages of the quorum include
2 validation timestamps that match.
- 1 6. The method of claim 1 wherein the stripe of erasure coded data was previously
2 stored using a technique of striping.
- 1 7. The method of claim 1 wherein the stripe of erasure coded data was not
2 previously stored using a technique of striping.
- 1 8. A method of reading data comprising the steps of:

2 receiving a request for a data block stored within a stripe of erasure
3 coded data, the stripe of erasure coded data stored across a plurality of
4 storage devices which include a target storage device that holds the data
5 block, the stripe comprising stripe blocks;

6 sending read messages to the storage devices, the read message sent to
7 the target storage device including an indicator; and

8 receiving a reply message from each of at least a quorum of the storage
9 devices which indicate that there is no pending write for the stripe block
10 stored on the storage device, the reply messages including validation
11 timestamps which match, the reply message from the target storage device
12 including the data block, the quorum meeting a quorum condition of a
13 number such that any two selections of the number of stripe blocks
14 intersect in the minimum number of the stripe blocks needed to decode the
15 stripe.

1 9. A method of writing data comprising the steps of:

2 receiving a data block for storage within a stripe of erasure coded data,
3 the stripe comprising stripe blocks;

4 sending a query message to each of a plurality of storage devices upon
5 which the stripe of erasure coded data is stored;

6 receiving a query reply message from each of at least a first quorum of
7 the storage devices;

8 sending a modify message to each of the storage devices; and

9 receiving a write reply message from each of at least a second quorum
10 of the storage devices, the first and second quorums each meeting a
11 quorum condition of a number such that any two selections of the number
12 of the stripe blocks intersect in a minimum number of the stripe blocks
13 needed to decode the stripe.

1 10. The method of claim 9 wherein the stripe blocks comprise a first number of
2 data blocks and a second number of parity blocks.

1 11. The method of claim 9 wherein the query message sent to the storage device
2 upon which the data block is to be stored includes an indicator.

- 1 12. The method of claim 9 wherein the query messages include a timestamp
2 indicating a current time.
- 1 13. The method of claim 12 wherein the reply messages of the first quorum each
2 indicate that the timestamp is later than a pending write timestamp.
- 1 14. The method of claim 12 wherein the reply messages of the first quorum each
2 indicate that the timestamp is later than a validation timestamp for a previous
3 version of the data block.
- 1 15. The method of claim 14 wherein the query reply message from the storage
2 device upon which the data block is to be stored includes the validation timestamp
3 and the previous version of the data block.
- 1 16. The method of claim 15 wherein each of the modify messages include the
2 timestamp and the validation timestamp.
- 1 17. The method of claim 16 wherein the modify message sent to the storage
2 device upon which the data block is to be stored includes the data block.
- 1 18. The method of claim 17 wherein the modify messages sent to the storage
2 devices which hold parity blocks include the previous version of the data block
3 and the data block.
- 1 19. The method of claim 18 wherein the write reply messages from the second
2 quorum indicate that the validation timestamp equals a maximum timestamp for
3 the stripe block stored on the storage device.
- 1 20. The method of claim 18 wherein the write reply messages from the second
2 quorum indicate that the timestamp is no earlier than the pending write timestamp.
- 1 21. The method of claim 18 wherein the write reply message from the storage
2 device which stores the data block indicates that the data block was stored

3 successfully.

1 22. The method of claim 17 wherein the modify messages sent to the storage
2 devices which hold parity blocks include a coded block which represents the
3 previous version of the data block and the data block.

1 23. The method of claim 9 wherein the stripe of erasure coded data was previously
2 stored using a technique of striping.

1 24. The method of claim 9 wherein the stripe of erasure coded data was not
2 previously stored using a technique of striping.

1 25. A method of writing data comprising the steps of:
2 receiving a particular data block for storage within a stripe of erasure
3 coded data, the stripe comprising stripe blocks which comprise a first
4 number of data blocks and a second number of parity blocks;
5 sending a query message including a timestamp indicating a current
6 time to each of a plurality of storage devices upon which the stripe of
7 erasure coded data is stored, the query message sent to a target storage
8 device upon which the particular data block is to be stored including an
9 indicator;
10 receiving a query reply message from each of at least a first quorum of
11 the storage devices indicating that the timestamp is later than a pending
12 write timestamp and that the timestamp is later than a validation timestamp
13 for the earlier version of the stripe block, the query reply message from the
14 target storage device including the validation timestamp and a previous
15 version of the particular data block;
16 sending a modify message to each of the storage devices including the
17 timestamp and the validation timestamp, the modify message sent to the
18 storage device upon which the block of data is to be stored including the
19 particular data block, the modify messages sent to the storage devices
20 which hold the parity blocks including information used to update the
21 parity blocks; and
22 receiving a write reply message from each of at least a second quorum

23 of the storage devices indicating that the validation timestamp equals a
24 maximum timestamp for the stripe block stored on the storage device, and
25 that the timestamp is no earlier than the pending write timestamp, the first
26 and second quorums each meeting a quorum condition of a number such
27 that any two selections of the number of the stripe blocks intersect in a
28 minimum number of the stripe blocks needed to decode the stripe.

1 26. The method of claim 25 wherein the information used to update the parity
2 blocks comprises a coded block which represents the previous version of the
3 particular data block and the particular data block.

1 27. The method of claim 25 wherein the information used to update the parity
2 blocks comprises the previous version of the particular data block and the
3 particular data block.

1 28. A computer readable memory comprising computer code for implementing a
2 method of reading data, the method of reading the data comprising the steps of:
3 receiving a request for a data block stored within a stripe of erasure
4 coded data, the stripe of erasure coded data stored across a plurality of
5 storage devices which include a target storage device that holds the data
6 block;
7 sending read messages to the storage devices; and
8 receiving at least a quorum of reply messages from the storage devices,
9 the reply message from the target storage device including the data block,
10 the quorum meeting a quorum condition of a number such that any two
11 selections of the number of stripe blocks intersect in a minimum number of
12 the stripe blocks needed to decode the stripe.

1 29. The computer readable memory of claim 28 wherein the stripe blocks
2 comprise a first number of data blocks and a second number of parity blocks.

1 30. The computer readable memory of claim 28 wherein the read message that is
2 sent to the target storage device includes an indicator.

1 31. The computer readable memory of claim 28 wherein the reply messages of the
2 quorum indicate that there is no pending write for the stripe block stored on the
3 storage device.

1 32. The computer readable memory of claim 28 wherein the reply messages of the
2 quorum include validation timestamps that match.

1 33. The computer readable memory of claim 28 wherein the stripe of erasure
2 coded data was previously stored using a technique of striping.

1 34. The computer readable memory of claim 28 wherein the stripe of erasure
2 coded data was not previously stored using a technique of striping.

1 35. A computer readable memory comprising computer code for implementing a
2 method of writing data, the method of writing the data comprising the steps of:
3 receiving a data block for storage within a stripe of erasure coded data,
4 the stripe comprising stripe blocks;
5 sending a query message to each of a plurality of storage devices upon
6 which the stripe of erasure coded data is stored;
7 receiving a query reply message from each of at least a first quorum of
8 the storage devices;
9 sending a modify message to each of the storage; and
10 receiving a write reply message from each of at least a second quorum
11 of the storage devices, the first and second quorums each meeting a
12 quorum condition of a number such that any two selections of the number
13 of the stripe blocks intersect in a minimum number of the stripe blocks
14 needed to decode the stripe.

1 36. The computer readable memory of claim 35 wherein the stripe blocks
2 comprise a first number of data blocks and a second number of parity blocks.

1 37. The computer readable memory of claim 35 wherein the query message sent to
2 the storage device upon which the data block is to be stored includes an indicator.

- 1 38. The computer readable memory of claim 35 wherein the query messages
2 include a timestamp indicating a current time.
- 1 39. The computer readable memory of claim 38 wherein the reply messages of the
2 first quorum each indicate that the timestamp is later than a pending write
3 timestamp.
- 1 40. The computer readable memory of claim 38 wherein the reply messages of the
2 first quorum each indicate that the timestamp is later than a validation timestamp
3 for a previous version of the data block.
- 1 41. The computer readable memory of claim 40 wherein the query reply message
2 from the storage device upon which the data block is to be stored includes the
3 validation timestamp and the previous version of the data block.
- 1 42. The computer readable memory of claim 41 wherein each of the modify
2 messages include the timestamp and the validation timestamp.
- 1 43. The computer readable memory of claim 42 wherein the modify message sent
2 to the storage device upon which the data block is to be stored includes the data
3 block.
- 1 44. The computer readable memory of claim 43 wherein the modify messages sent
2 to the storage devices which hold parity blocks include the previous version of the
3 data block and the data block.
- 1 45. The computer readable memory of claim 44 wherein the write reply messages
2 from the second quorum indicate that the validation timestamp equals a maximum
3 timestamp for the stripe block stored on the storage device.
- 1 46. The computer readable memory of claim 44 wherein the write reply messages
2 from the second quorum indicate that the timestamp is no earlier than the pending
3 write timestamp.

1 47. The computer readable memory of claim 44 wherein the write reply message
2 from the storage device which stores the data block indicates that the data block
3 was stored successfully.

1 48. The computer readable memory of claim 43 wherein the modify messages sent
2 to the storage devices which hold parity blocks include a coded block which
3 represents the previous version of the data block and the data block.

1 49. The computer readable memory of claim 35 wherein the stripe of erasure
2 coded data was previously stored using a technique of striping.

1 50. The computer readable memory of claim 35 wherein the stripe of erasure
2 coded data was not previously stored using a technique of striping.